

# 1 **Single-cell Long Non-coding RNA Landscape of T Cells in Human** 2 **Cancer Immunity**

3 Haitao Luo<sup>1,3,4,\*,#a</sup>, Dechao Bu<sup>2,#b</sup>, Lijuan Shao<sup>1,3,4,#c</sup>, Yang Li<sup>5,d</sup>, Liang Sun<sup>2,e</sup>, Ce Wang<sup>1,3,f</sup>,  
4 Jing Wang<sup>1,3,4,g</sup>, Wei Yang<sup>1,3,h</sup>, Xiaofei Yang<sup>1,3,i</sup>, Jun Dong<sup>4,\*j</sup>, Yi Zhao<sup>2,\*k</sup> and Furong Li<sup>1,3,\*l</sup>

5 <sup>1</sup> *Translational Medicine Collaborative Innovation Center, The Second Clinical Medical*  
6 *College (Shenzhen People's Hospital), Jinan University, Shenzhen 518020, China*

7 <sup>2</sup> *Bioinformatics Research Group, Key Laboratory of Intelligent Information Processing,*  
8 *Advanced Computing Research Center, Institute of Computing Technology, Chinese*  
9 *Academy of Sciences, Beijing 100190, China*

10 <sup>3</sup> *Shenzhen key laboratory of stem cell research and clinical transformation, Shenzhen*  
11 *518020, China*

12 <sup>4</sup> *Integrated Chinese and Western Medicine Postdoctoral research station, Jinan University,*  
13 *Guangzhou 510632, China*

14 <sup>5</sup> *Department of Gastrointestinal Surgery, The Second Clinical Medical College (Shenzhen*  
15 *People's Hospital), Jinan University, Shenzhen 518020, China*

16

17 \* Corresponding authors.

18 E-mail: [luoht1985@gmail.com](mailto:luoht1985@gmail.com) (Luo H), [dongjunbox@163.com](mailto:dongjunbox@163.com) (Dong J), [biozy@ict.ac.cn](mailto:biozy@ict.ac.cn)  
19 (Zhao Y), [frli62@163.com](mailto:frli62@163.com) (Li F).

20 # Equal contribution.

21 **Running title:** *Luo H et al / Single-cell LncRNA Landscape of T Cells*

22

23 Total word counts (from “Introduction” to “Materials and methods”): 5,625

24 Total figures: 6

25 Total tables: 0

26 Total supplementary figures: 6

27 Total supplementary tables: 11

28 **Abstract**

29 The development of new therapeutic targets for cancer immunotherapies and the  
30 development of new biomarkers require deep understanding of T cells. To date, the complete  
31 landscape and systematic characterization of long noncoding RNAs (lncRNAs) in T cells in  
32 cancer immunity are lacking. Here, by systematically analyzing full-length single-cell RNA  
33 sequencing (scRNA-seq) data of more than 20,000 T cell libraries across three cancer types,  
34 we provide the first comprehensive catalog and the functional repertoires of lncRNAs in  
35 human T cells. Specifically, we developed a custom pipeline for *de novo* transcriptome  
36 assembly obtaining 9,433 novel lncRNA genes that increased the number of current human  
37 lncRNA catalog by 16% and nearly doubled the number of lncRNAs expressed in T cells. We  
38 found that a portion of expressed genes in single T cells were lncRNAs which have been  
39 overlooked by the majority of previous studies. Based on metacell maps constructed by  
40 MetaCell algorithm that partition scRNA-seq datasets into disjointed and homogenous groups  
41 of cells (metacells), 154 signature lncRNAs associated with effector, exhausted, and  
42 regulatory T cell states are identified, 84 of which are functionally annotated based on co-  
43 expression network, indicating that lncRNAs might broadly participate in regulation of T cell  
44 functions. Our findings provide a new point of view and resource for investigating the  
45 mechanisms of T cell regulation in cancer immunity as well as for novel cancer-immune  
46 biomarker development and cancer immunotherapies.

47 **KEYWORDS:** lncRNA; Transcriptome assembly; Metacell; Immune regulation; Functional  
48 annotation

49

## 50 **Introduction**

51 T cell checkpoint inhibition therapies, such as targeting exhausted CD8<sup>+</sup> T cells and  
52 regulatory T cells (Tregs), have shown remarkable clinical benefit in many cancers [1-3].  
53 Nevertheless, the mechanisms underlying therapy response or resistance are largely unknown,  
54 which leads to the different therapeutic efficacies among cancer patients [4-8]. To better  
55 understand the mechanisms that underlie successful response to immunotherapy, more  
56 comprehensive studies to explore the whole transcriptome of individual T cells in tumor  
57 ecosystems are desired. Long non-coding RNAs (lncRNAs), defined as a class of non-coding  
58 RNAs longer than 200 nucleotides with no or low protein-coding potential, comprise a large  
59 proportion of the mammalian transcriptome [9-12]. Accumulating evidence has suggested  
60 that lncRNAs are widely expressed in immune cells and play crucial roles in cancer immunity  
61 by regulating the differentiation and function of T cells [13-17]. For example, overexpression  
62 of *NKILA*, an *NF-κB*-interacting lncRNA, correlated with T cell apoptosis and shorter patient  
63 survival [18], and an enhancer-like lncRNA *NeST* regulates epigenetic marking patterns of  
64 *IFN-γ*-encoding chromatin and induce synthesis of *IFN-γ* in CD8 T cells [19]. However,  
65 previous studies seem to be somewhat scattered and the landscape and comprehensive  
66 functional analysis of lncRNAs in T cells in cancer immunity are still lacking.

67 The dramatic advances of single-cell RNA sequencing (scRNA-seq) technologies have  
68 gained unprecedented insight into the high diversity in T cell types and states compared to  
69 bulk RNA sequencing methods, which do not address the complex structures of tumor  
70 microenvironment [20-25]. Despite the advantages of single-cell resolution, in current most  
71 scRNA-seq studies of cancer immunology have generally focused on coding genes,  
72 overlooking the large amounts of lncRNAs. Detailed understanding of lncRNAs at the single-  
73 cell level was challenging owing to their relatively low and cell-specific expression [26-28].  
74 As a widely used scRNA-seq approach, 3'-end sequencing technologies such as droplet-  
75 based 10X Genomics have lower RNA capture efficiencies, leading to the dropout events and  
76 technological noise for lowly expressed lncRNAs [29]. Furthermore, accurate identification  
77 of novel lncRNAs is not suitable for the 3'-end sequencing technologies, but such analysis  
78 could be achieved by using full-length scRNA-seq technologies such as SMART-seq2 [30].  
79 In addition, the sampling noise in scRNA-seq is generated through sampling of limited RNA  
80 transcripts from each cell [31], leading to a highly noisy estimation for most lncRNAs.  
81 Therefore, to effectively characterize the lncRNA landscape at single-cell level, attention  
82 should be paid to choosing appropriately scRNA-seq data and analytical approaches.

83 Here, using unprecedentedly large-scale full-length single-cell transcriptome data of more  
84 than 20,000 T cells from various tissues across three cancer types, we created a full  
85 annotation of the T cell lncRNA transcriptome and analyzed the functional roles associated  
86 with different T cell states. Our study aims to provide a basic and valuable resource for the  
87 future exploration of lncRNA regulatory mechanisms in T cells, which may facilitate novel  
88 cancer-immune biomarker development.

89

## 90 **Results**

### 91 ***De novo* transcriptome assembly of lncRNAs from scRNA-seq data of T cells**

92 To investigate the landscape of human lncRNAs in T cells across different tissues, patients  
93 and cancer types, we collected the data of 24,068 T cells (the size of the gzip-compressed  
94 FASTQ file was 7.5 TB) generated by full-length single-cell RNA sequencing with SMART-  
95 seq2, including the raw data of 9,878 cells from colorectal cancer (CRC) patients (2.8 TB) ,  
96 10,188 cells from non-small-cell lung cancer (NSCLC) patients (3.1 TB), and 4,002 cells  
97 from 5 hepatocellular carcinoma (HCC) patients (1.6 TB) [32-34] (Figure S1A and Table S1).  
98 These cells were collected from peripheral blood, adjacent normal, and tumor tissue from  
99 each patient and sorted into CD3<sup>+</sup>CD8<sup>+</sup> (CD8) and CD3<sup>+</sup>CD4<sup>+</sup> (CD4) T cells. The reads of  
100 each cell were mapped to the human reference genome (hg38/GRCh38), and the cells with  
101 unique mapping rates of less than 20% were removed. The remaining cells with on average  
102 1.04 million uniquely mapped read pairs (0.63 million splices on average) and at least one  
103 pair of T cell receptor (TCR)  $\alpha$  and  $\beta$  chains enabled us to detect the expressed lncRNAs  
104 (Figure S1B-D).

105 Next, to generate the comprehensive T cell transcriptome beyond the currently reference  
106 annotation, we performed *de novo* transcriptome assembly using the StringTie method [35].  
107 Although StringTie could be run by providing the reference annotation to guide the transcript  
108 construction, in current study we focused on to what extent it could assemble the whole  
109 transcriptome without the prior annotation. Based on the T cell dataset from HCC patients,  
110 we first measured the extent of assembly in each T cell and found that an average of 4,752  
111 transcripts could be assembled at single-cell level, and an average of 69.8% (3,318/4,752)  
112 were matched to reference models (including reference protein-coding genes from  
113 GENCODE v31 and reference lncRNA genes from RefLnc database) (**Figure 1A**).

114 To explore the best way to obtain novel transcripts, we compared the assembly results  
115 using three different approaches based on HCC dataset: (1) mapping and assembling for each  
116 single cell individually (cell-level); (2) assembling transcripts based on merged mapping  
117 results from each cell type of each patient (cell type-level); (3) assembling transcripts based  
118 on merged mapping results from each tissue of each patient (tissue-level). The transcripts  
119 assembled from each approach were merged independently and compared with reference  
120 genes respectively (Figure S1E). We found that the number of assembled transcripts  
121 matching to reference genes based on the cell type-level strategy (average 105,527 transcripts)  
122 was significantly higher than in cell-level or tissue-level methods (average 77,860 and 49,689  
123 transcripts respectively,  $P$ -value  $< 0.001$ , Wilcoxon rank sum test) (Figure 1B). Furthermore,  
124 the average number of matched transcripts from the cell type-level was more than twice that  
125 from the bulk-seq method (average 48,854 transcripts) (Figure 1B).

126 According to the cell-type pooling strategy, the cells from all patients across three cancer  
127 types were partitioned into 266 subsets (Figure 1C and Figure S1A), and the mapping results  
128 of cells from the same subset were merged and fed into assembling program. We found the  
129 number of assembled transcripts across different subsets showed positive correlations with  
130 the number of cells in these subsets in both CRC and NSCLC datasets (Pearson correlation  
131 coefficients = 0.6 and 0.72,  $P$ -value =  $4.3e-11$  and  $< 2.2e-16$ , respectively), but not in the  
132 HCC dataset (Pearson correlation coefficient = 0.22,  $P$ -value = 0.17) (Figure 1D and Figure  
133 S1F). Then, assembled transcripts from all subsets were merged together, and a total of  
134 751,710 primary genes were obtained. Next, we compared our assembled transcriptome with  
135 reference gene models. The results showed that reference lncRNAs had a lower detection rate  
136 than protein-coding genes. Specifically, 82% (16,399/19,938) of the known protein-coding  
137 genes in GENCODE v31 could be verified (44%, 8,893/19,938 were complete match with the  
138 same intron chain), while 16% (9,567/59,489) of known lncRNA genes were verified (5%,  
139 3,140/59,489 were complete match) (Figure 1E). These findings suggested that lncRNAs  
140 were expressed in a much more cell-specific manner than protein-coding genes and further  
141 studies to uncover novel lncRNAs specifically expressed in human T cells were needed.

142 From the primary assembly, we developed a custom pipeline to identify novel lncRNAs.  
143 Briefly, we first selected transcripts that were no shorter than 200 nucleotides and have  
144 multiple exons. The transcripts that overlapped with both known protein-coding and known  
145 lncRNA genes were filtered out. Then, the transcripts lacking coding potential predicted by  
146 both CPC [36] and CNCI [37] utility were retained. Finally, the remaining transcripts that

147 were reconstructed in at least two subsets with complete match were defined as the novel  
148 lncRNA catalog (Figure 1C). Through this multi-layered analysis, we identified 9,433  
149 previously unknown lncRNA genes (13,025 transcripts with mean length of 1,112  
150 nucleotides), which increased the number of current human lncRNA catalog [38] by 16% and  
151 nearly doubled the number of lncRNAs expressed in human T cells.

152 Finally, we performed experimental validation to evaluate the robustness of our identified  
153 novel lncRNAs. First, fresh peripheral blood samples were collected from three CRC patients  
154 (Table S2). Then, mononuclear cells were isolated from each sample. CD8 and CD4 T cells  
155 were separated by immunomagnetic beads and the separation efficiency was verified by flow  
156 cytometry (Figure 2A). Next, we selected 50 novel lncRNAs for quantitative real-time  
157 polymerase chain reaction (qRT-PCR) analysis and Sanger sequencing across T cell samples.  
158 As a result, 38 novel lncRNAs could be verified successfully by Sanger sequencing (Table  
159 S3). As an example, for a novel lncRNA *TCONS\_00180551* located in an intergenic region of  
160 chromosome 11, the blat search result of Sanger sequencing exactly matches the junction of  
161 this novel lncRNA (Figure 2B).

## 162 **The characterization and expression analyses of lncRNAs in T cells**

163 Based on the relative genomic locations to reference protein-coding genes, the novel  
164 lncRNAs were classified into three locus biotypes, including 6,525 as intergenic, 3,187 as  
165 intronic and 3,313 as antisense lncRNAs. As in the case of reference lncRNAs, these novel  
166 lncRNAs showed fewer exons (the average number of exons was 2), lower detection rates  
167 and average gene abundance than protein-coding genes at single-cell level (Figure 3A-B).  
168 Specifically, by using pseudoalignment of scRNA-seq reads to both reference and novel  
169 lncRNA transcriptomes, on average 5,902 genes were detected (counts larger than 1) in each  
170 cell, 41% (2,397) of which were lncRNAs, including 1,258 reference and 1,139 novel  
171 lncRNAs (Figure 3A). Furthermore, for both reference and novel lncRNA genes, the average  
172 number of expressed genes across T cells was significantly lower than that of protein-coding  
173 genes. More precisely, we found that an average of 5,596 protein-coding and 2,093 lncRNA  
174 genes were expressed in at least 25% of cells. In such a situation, novel lncRNAs exhibited a  
175 higher average expression number and expression rate (1,489 and 15.8%, 1,489/9,433) than  
176 did reference lncRNAs (604 and 1%, 604/59,489) (Figure 3B), suggesting that novel  
177 lncRNAs exhibited more enrichment than known lncRNAs in T cells in cancer. Moreover,  
178 we performed further analysis to investigate the specifically expressed lncRNAs in different

179 tissues for each cancer type. In brief, for both CD4 and CD8 T cells of each cancer type, we  
180 identified 96 and 90 lncRNAs including 44 and 40 novel lncRNAs that expressed in tissue-  
181 specific pattern (Table S4). For example, some novel lncRNAs such as *XLOC-301694* and  
182 *XLOC-126527* were significantly expressed in CD4 T cells from tumor tissue of CRC  
183 (adjusted *P* value =3.17E-68 and 1.72E-64 respectively), while others such as *XLOC-302096*  
184 and *XLOC-502999* were significantly enriched in normal tissue and peripheral blood  
185 respectively (adjusted *P* value =9.18E-82 and 1.35E-44 respectively) (Table S4). Finally, we  
186 assessed the evolutionary conservation of these novel lncRNA transcripts and found that, on  
187 average, 61.2% have orthologous regions in the primate genomes, while only 3.4% mapped  
188 to mouse genome, suggesting the poor sequence conservation of these novel lncRNAs.

### 189 **Identification of signature lncRNAs associated with T cell states in cancer immunity** 190 **based on metacell maps**

191 To explore signature lncRNAs associated with T cell states in cancer immunity, we used the  
192 MetaCell method [31] that partitioned the scRNA-seq datasets into disjointed and  
193 homogeneous cell groups (namely metacells) using the non-parametric *K*-nn graph algorithm.  
194 For the lowly and specifically expressed nature of lncRNA genes, metacells pooling together  
195 data from cells derived from the same transcriptional states could serve as building blocks for  
196 approximating the distributions of lncRNA gene expression and minimizing the technical  
197 variance and noise. After quality control, 19,572 cells with predefined cluster annotations and  
198 21,205 genes including both protein-coding and lncRNA genes were retained and used for the  
199 following analyses. The expression tables of CD8 and CD4 T cells across three cancers were  
200 fed into the MetaCell pipeline separately, resulting in a detailed map of 43 and 65 metacells  
201 respectively (**Figure 4A-B** and Table S5-6).

202 Based on the 2D projections (Figure 4A-B), predefined cell cluster annotations (Table  
203 S1), and the metacell similarity matrices (similarity among 43 or 65 metacells for CD8 or  
204 CD4 T cells respectively) (Figure S2A-B and Figure 4C-D), we organized the complex  
205 transcriptional landscape of CD8 into Naïve, effector/pre-effector, intermediated, and  
206 exhausted metacell groups and CD4 into Naïve, effector, intermediated, exhausted, and  
207 regulatory (including *FOXP3*<sup>+</sup>*CTLA4*<sup>low</sup> and *FOXP3*<sup>+</sup>*CTLA4*<sup>high</sup>) metacell groups respectively  
208 (Figure 4C-D). To evaluate the composition of metacells, we mapped tissue- and cancer-  
209 specific patterns in all metacells and achieved results in accordance with previous studies  
210 [32-34] (Figure 4C-D and Figure S3-4). For example, exhausted metacells were preferentially

211 enriched in tumors, while effector metacells were prevalent in peripheral blood. Although  
212 some metacells were enriched in different cancer types, they were organized into the same  
213 functional groups (Figure 4C-D). Notably, effector metacell groups (cytotoxic state) and  
214 exhausted metacell groups (dysfunctional state) were located in different directions in the  
215 metacell maps, while the diffuse border was observed between the intermediate and the  
216 cytotoxic or dysfunctional state (Figure 4E-F). These intermediate cells exhibited remarkable  
217 transcriptional heterogeneity indicating functional divergence of these cells (Figure 4E-F and  
218 Figure S3-4). The observed cluster distribution in both CD8 and CD4 metacell maps might  
219 suggest a relative transition from activation to exhaustion that began with Naïve cells,  
220 followed by intermediate cells (such as central memory (CM), effector memory (EM) and  
221 tissue resident memory (RM) cells) and ended with exhausted cells. Moreover, the CD4  
222 metacell map revealed that Tregs were subdivided into *FOXP3*<sup>+</sup>*CTLA4*<sup>low</sup> Tregs and  
223 *FOXP3*<sup>+</sup>*CTLA4*<sup>high</sup> Tregs that were preferentially enriched in blood and tumors respectively  
224 (Figure 4D and 4F). These observations demonstrated that the diversity and dynamics of T  
225 cell states in cancer immune infiltrates could be controlled by complex and intricate gene  
226 regulatory mechanisms. Yet, the association between these cell states and lncRNAs was still  
227 poorly characterized, prompting us to subsequently investigate potential roles of lncRNA  
228 genes in T cells. Currently, the cell groups such as *FOXP3*<sup>+</sup>*CTLA4*<sup>high</sup> Tregs and exhausted T  
229 cells expressing inhibitory receptors (e.g., *PDCD1* and *TIGIT*) have been used as therapeutic  
230 targets for anti-cancer immunotherapies, thus we focused on these cells in the following  
231 analyses.

232 To explore signature lncRNAs associated with effector T cells, exhausted T cells, and  
233 Tregs, we performed systematic analysis of these metacell groups based on well-defined  
234 anchor genes [39], such as the genes associated with CD8 effector functions (*CX3CRI*,  
235 *FGFBP2*, *GZMH* and *PRF1*) or with the CD8 exhausted state (*HAVCR2*, *LAG3*, *PDCD1*,  
236 *TIGIT* and *CTLA4*). As a result, 154 lncRNAs that were significantly correlated to the anchor  
237 genes were identified and were involved in a set of co-expressed gene modules, including  
238 effector, exhausted and Treg gene modules (**Figure 5A-B** and Table S7). Interestingly, a  
239 putative *CTLA4*<sup>high</sup> Treg gene subset was observed in the Treg module, suggesting its specific  
240 functional roles in tumor-infiltrating Treg cells (Figure 5B). Overall, by combination analysis  
241 with the expression profile across metacell groups, we found 47 and 79 lncRNAs correlated  
242 with effector and exhausted states in CD8 and CD4 cells respectively and were designated as  
243 effector or exhausted signature lncRNAs (Figure 5C and Figure S5). Similarly, 49 lncRNAs



244 were highly associated with Treg cells and were designated as Treg signature lncRNAs  
245 (Figure S5). Among these signature lncRNAs, 14 and 7 lncRNA genes were shared between  
246 CD8 and CD4 effector states and between CD8 and CD4 exhausted states respectively. 21  
247 lncRNA genes associated with Tregs overlapped with those characteristics in the exhausted  
248 CD4 T cells (Table S7), indicating the presence of shared regulatory roles of these lncRNAs.  
249 In contrast, no signature lncRNA was shared between exhausted and effector states.

### 250 **Functional prediction of signature lncRNAs associated with T cell states based on co-** 251 **expression network**

252 To gain further insights into the functional roles of lncRNA in different T cell states in cancer,  
253 we built a coding-noncoding network (CNC), as we previously reported [40, 41], using linear  
254 correlation over all metacells. Applying this strategy, the functions of 54% (84/154) signature  
255 lncRNAs were annotated (Table S8). As expected, both CD8 and CD4 exhausted T cells have  
256 the functional enrichments of signature lncRNAs that were markedly different from effector  
257 CD8 or CD4 T cells, including regulation manners in immune system processes and several  
258 signalling pathways (**Figure 6A-B**). For example, exhausted signature lncRNAs were  
259 significantly enriched in immunoinhibitory functions such as negative regulation of immune  
260 response (adjusted  $P$ -value =  $2.96e-14$ ), negative regulation of T cell activation (adjusted  $P$ -  
261 value =  $1.24e-06$ ), and positive regulation of interleukin-10 biosynthetic process (adjusted  $P$ -  
262 value =  $1.02e-18$ ). In comparison, effector signature lncRNAs were enriched in cytotoxic  
263 programs such as T cell proliferation involved in immune response (adjusted  $P$ -value =  
264  $8.16e-09$ ), positive regulation of cytokine secretion (adjusted  $P$ -value =  $4.65e-05$ ), and  
265 positive regulation of cytolysis (adjusted  $P$ -value =  $1.59e-19$ ) (Figure 6A-B and Table S9-10).  
266 These results consisted with the phenotypes of exhausted or effector states of T cells as  
267 described in previous studies [1, 32-34, 42]. In addition, the enriched functions of Treg  
268 signature lncRNAs were similar with those of CD4 exhausted signature lncRNAs involving  
269 multiple immunosuppressive programs (**Figure 6C** and Table S11), suggesting the shared  
270 regulatory roles of these lncRNAs in CD4 Tregs and exhausted CD4 T cells. Further analysis  
271 of the functions of co-signature lncRNAs that were shared between CD8 and CD4 exhausted  
272 or effector states, as well as between CD4 exhausted and Treg states (Figure S6), suggests  
273 that the signature lncRNAs might broadly participate in regulation of T cell functions within  
274 the human tumor microenvironment.

275 For example, a known lncRNA *TM4SF19-ASI*, defined as a signature lncRNA for both  
276 CD8 effector and CD4 effector T cells and was transcribed in the antisense orientation to the  
277 *TM4SF19* gene, was co-expressed with 66 protein-coding and 11 lncRNA genes (Figure 6D-  
278 E). Of note, *TM4SF19-ASI* was highly correlated and located in the same topologically  
279 associated domain (TAD) with its host gene *TM4SF19* (Pearson correlation coefficient = 0.88)  
280 (Figure 6D), a member of the four-transmembrane L6 superfamily participating in various  
281 cellular processes including cell proliferation, motility, and cell adhesion [43-46].  
282 Consistently, *TM4SF19-ASI* was significantly enriched in several effector T cell associated  
283 processes such as cellular response to cholesterol (adjusted *P*-value = 1.09e-30), cell adhesion  
284 (adjusted *P*-value = 5.25e-27) and regulation of tumor necrosis factor biosynthetic process  
285 (adjusted *P*-value = 3.75e-11) (Figure 6F). Interestingly, a recent study suggested that anti-  
286 tumor response of CD8 T cells could be enhanced by regulating cholesterol metabolism [47].  
287 For another example, a novel lncRNA *XLOC-633950*, defined as a signature lncRNA for  
288 both CD4 exhausted T cells and Treg cells, was an intergenic gene and transcribed from the  
289 promoter-enhancer cluster region of the *SLA* and *CCN4* genes (Figure 6G). Furthermore,  
290 *XLOC-633950* as a novel gene, whose expression was supported by multiple expressed  
291 sequence tags (EST), was located in the same TAD with the *SLA* gene which acted as an  
292 inhibitor of antigen receptor signalling by negative regulation of positive selection and  
293 mitosis of T cells [48-51] (Figure 6G). In accordance with *SLA* functions, the functional  
294 enrichments of *XLOC-633950* according to its co-expressed protein-coding genes were  
295 mainly associated with immunoinhibitory processes, such as negative regulation of T cell  
296 cytokine production (adjusted *P*-value = 4.56e-13) and negative regulation of T cell  
297 proliferation and activation (adjusted *P*-value = 7.25e-11 and 5.85e-08 respectively) (Figure  
298 6H-I). These results provided a starting point for future dissecting the mechanisms of  
299 signature lncRNAs.

## 300 Discussion

301 Despite the obvious advantages, most scRNA-seq data was still limited in its ability to study  
302 lncRNAs, which were emerging as central players and key regulators in a number of  
303 biological processes such as anti-tumor immune response [52, 53]. In comparison with many  
304 scRNA-seq methods that amplified only the 3' end of transcripts, the SMART-seq2 protocol  
305 could generate full-length cDNA from polyadenylated transcripts which results in data  
306 suitable for analysis of lncRNAs [30, 54]. In current study, we preformed systematic analyses

307 of SMART-seq2 full-length scRNA-seq datasets and provided the first comprehensive atlas  
308 of lncRNA in T cells of human cancer.

309 Recently, Jiang *et al.* presented a comprehensive human lncRNA catalog (RefLnc) [38]  
310 containing 77,900 lncRNAs based on analysis of 14,166 polyA(+) RNA-Seq libraries and  
311 previous known annotations. Among the RefLnc lncRNAs, only 16% could be assembled and  
312 expressed in T cells. In addition, compared with bulk-seq data, scRNA-seq data could  
313 detect more known and novel transcripts. These observations suggested that despite the  
314 vast number of lncRNAs that have been identified using bulk-seq data [10, 12, 26, 38, 55],  
315 the catalog of human lncRNAs is still far from being complete at single-cell resolution, due to  
316 their low and cell-specific expression patterns. Based on the cell-pooling strategy and more  
317 than 20,000 scRNA-seq libraries from 31 patients across three cancer types, we identified  
318 9,433 previously non-annotated lncRNAs. These results significantly expand the current  
319 lncRNA catalog and enable us to carry out in-deep analysis of the T cell context-specific  
320 lncRNA transcriptome. Notably, all the scRNA-seq data used in current study was generated  
321 by sequencing the polyadenylated (polyA) transcriptome, in which non polyadenylated  
322 lncRNAs were absent.

323 Several previous studies have applied full-length scRNA-seq to unleash tumor infiltrating  
324 lymphocytes in HCC [34], NSCLC [32], and CRC[33], providing a deep understanding of the  
325 immune landscape of T cells in cancer. Nevertheless, the physiological function of lncRNAs  
326 in different T cell states during the cancer immune response remains elusive. Although the  
327 abundance of lncRNA was relatively low and hard to distinguish from technical noise in  
328 single T cells, pooling the transcripts from multiple cells that are derived from the same cell  
329 state allows more accurate quantification of lncRNAs, making it feasible to explore their  
330 signatures and putative regulatory mechanisms associated with T cell states in cancer  
331 immunity. Based on such partitioning and pooling strategies, we used the MetaCell method to  
332 identify homogeneous T cell groups from scRNA-seq data and derived a detailed map of 43  
333 and 65 metacells for CD8 and CD4 T cells respectively. These metacells with higher  
334 homogeneity, allowed a more accurate quantification of lncRNAs as well as identification of  
335 T cell differentiation gradients. For example, we observed 7 metacells involved in CD8  
336 effector cell cluster, which might reflect the transcriptional heterogeneity in this cluster  
337 (Figure 4C). The roles of lncRNAs in these different subsets (metacells) of CD8 effector T  
338 cells need further investigation. While MetaCell was not designed to perform single-cell

339 lncRNA analysis, the MetaCell partitioning algorithm facilitated robust cell grouping of  
340 scRNA-seq data which enabled us to study lncRNAs more accurately.

341 According to the metacell maps (Figure 4E-F), in contrast to the pool of intermediate T  
342 cells with diffuse borders with other cell states, a discrete pool of effector T cells, exhausted  
343 T cells and Tregs were observed that show clear gaps among them, thus facilitating unbiased  
344 analysis of signature lncRNAs in these cell states. In total, the 154 signature lncRNAs were  
345 obtained providing a useful reference lncRNA resource to further investigate their functions  
346 in T cell mediated cancer immunity. Since lncRNAs generally interact with protein-coding  
347 genes, and highly correlated genes generally have similar functions, the putative functions of  
348 these signature lncRNAs could be predicted by the co-expressed coding genes. Therefore, by  
349 constructing the ‘two color’ co-expression network in which both coding and lncRNA genes  
350 were involved, the functions of 84 signature lncRNAs were annotated. Some lncRNAs were  
351 genomically co-located with their host genes, that revealed the complicated regulation  
352 mechanisms of lncRNAs in cancer immunity. For example, as described above, *TM4SF19-*  
353 *ASI* was both co-expressed and co-located with their host gene *TM4SF19*, whose family has  
354 functions in various biological processes including cell proliferation and adhesion that are  
355 consistent with the characteristics of effector T cells [43-46].

356 In summary, the current study provides the first comprehensive catalog and the functional  
357 repertoires of lncRNAs in human cancer T cells. Although the expression pattern and exact  
358 mechanisms of these signature lncRNAs in regulating T cell states needs further experimental  
359 validation, we provide the groundwork for future studies to investigate the functional  
360 mechanisms of lncRNAs in the T cell mediated cancer immunity, especially in two of the  
361 essential states of T cells: effector state and exhausted state. These signature lncRNAs of  
362 CD8 exhausted T cells and tumor Tregs, may serve as new targets for novel cancer-immune  
363 biomarker development and cancer immunotherapies.

## 364 **Materials and methods**

### 365 **Full-length scRNA-seq and bulk RNA-seq datasets from cancer patients**

366 Raw sequencing data of three compendium datasets used in the current study were authorized  
367 by the European Genome-phenome Archive (EGA) and obtained from the EGA database  
368 under study accession id: EGAS00001002791, EGAS00001002430, and EGAS00001002072.  
369 The CRC scRNA-seq dataset (EGAS00001002791) contains the raw data of 11,138 single T

370 cells isolated from different tissues (peripheral blood, adjacent normal and tumor tissues) of  
371 12 CRC patients [33]. The NSCLC scRNA-seq dataset (EGAS00001002430) contains the  
372 raw data of 12,346 single T cells from 14 NSCLC patients [32]. The HCC scRNA-seq dataset  
373 (EGAS00001002072) contains the raw data of 5,063 single T cells from 6 HCC patients [34].  
374 All the data were generated by Illumina HiSeq 2500 sequencer with 100 bp pair-end reads or  
375 Illumina HiSeq 4000 sequencer with 150 bp pair-end reads. The cells from HCC patient  
376 P1202 (TCRs could not be assembled in those cells) were not analyzed in the current study.  
377 After preliminary filtration, 24,075 T cells with at least one pair of TCR *alpha-beta* chain  
378 were retained. The bulk RNA-seq data of five tumor samples from HCC patients were  
379 obtained from HCC dataset.

380 According to the cell annotations from original papers [32-34], these T cells were  
381 classified into different subtypes (Figure S1A and Table S1). PTC, NTC, and TTC represent  
382 CD3<sup>+</sup>CD8<sup>+</sup> T cells that were isolated from peripheral blood, adjacent normal, and tumor  
383 tissues respectively. The PTH, NTH, and TTH represent CD3<sup>+</sup>CD4<sup>+</sup>CD25<sup>low</sup> T cells that were  
384 isolated from the three tissues. PTR, NTR, and TTR represent CD3<sup>+</sup>CD4<sup>+</sup>CD25<sup>high</sup> T cells  
385 that were isolated from the three tissues.

### 386 **Reads mapping and transcripts assembly**

387 Clean reads from each T cell were mapped to the human reference genome (version  
388 hg38/GRCh38) using STAR aligner (v2.7.1) [56] with the *twopassMode* set as Basic. The  
389 bam files of T cells from each cell-type of each patient were merged using SAMtools merge  
390 [57]. StringTie (v2.0.3) [35] was used to assemble transcripts based on genomic read  
391 alignments. Assembled transcripts of all cell-types across all patients were merged together  
392 using the Cuffmerge utility of Cufflinks package [58].

### 393 **Comparison with reference gene annotation**

394 For reference gene annotation, lncRNA genes were collected from RefLnc [38] and other  
395 genes were collected from GENCODE v31 [59]. According to the “class code” information  
396 outputted by Cuffcompare, the merged assembly was classified into four categories by  
397 comparison with the reference gene annotation, including known coding genes, known  
398 lncRNA genes, potentially novel genes (class code is “i, x, u”), and others.

### 399 **Identification of novel lncRNAs**

400 Based on the potentially novel gene catalog derived from single-cell data, we developed a  
401 custom pipeline for identification of reliable novel lncRNAs including the following steps: (1)  
402 transcripts that are no shorter than 200 nt and have more than one exon were selected for  
403 downstream analysis (for intergenic transcripts, at least 1 kb away from known protein-  
404 coding genes); (2) CPC (Coding Potential Calculator) [36] and CNCI (Coding Noncoding  
405 Index) [37] software were used to evaluate the protein-coding potential of transcripts, and  
406 transcripts that were reported to lack coding potential by both CPC and CNCI were regarded  
407 as candidate noncoding transcripts; (3) The remaining transcripts that were assembled and  
408 have the same intron chain of at least two cell-types were retained as the final novel lncRNA  
409 catalog. The final lncRNA catalog was obtained by combining the reference lncRNA and  
410 novel lncRNA genes directly. The UCSC liftOver tool ([http://genome.ucsc.edu/cgi-  
411 bin/hgLiftOver?hgslid=806106955\\_h2xhcK2iPRI7SiMkxkB41I2mwF9O](http://genome.ucsc.edu/cgi-bin/hgLiftOver?hgslid=806106955_h2xhcK2iPRI7SiMkxkB41I2mwF9O)) was used to  
412 identify the orthologous locations of human novel lncRNAs in the mouse genome and in  
413 primates such as chimpanzee and gorilla, with the parameters: Minimum ratio of bases that  
414 must remap = 0.1 and Min ratio of alignment blocks or exons that must map =0.5.

#### 415 **Experimental validation of novel lncRNAs**

416 Three CRC patients were enrolled at Shenzhen People's Hospital. The informed consent  
417 forms were provided by patients. The current study was approved by Medical Ethics  
418 Committee of Shenzhen People's Hospital. The clinical characteristics of three patients are  
419 summarized in Table S2. Peripheral blood samples from three patients were obtained and  
420 treated with anticoagulation. Peripheral blood mononuclear cells (PBMCs) were extracted by  
421 Ficoll-Paque Plus (GE Healthcare, Sweden, 17144003). Then, CD8<sup>+</sup> and CD4<sup>+</sup> T cells were  
422 separated by immunomagnetic beads (Miltenyi Biotec, Germany, 130045101, 130045101).  
423 The separation efficiency was verified by flow cytometry. The sorted cells were dissolved in  
424 Trizol Reagent (Ambion, USA, 15596026) for RNA extraction according to the  
425 manufacture's protocol. cDNA was synthesized by PrimerScript RT reagent kit (Takara,  
426 Japan, AHG1552A). We chose 50 novel lncRNAs to perform experimental validation  
427 according to the following criteria: (1) highly expressed in either CD8 or CD4 T cells; (2)  
428 reconstructed in at least ten subsets with complete match; (3) uniquely mapped to human  
429 genome. For each lncRNA, at least two pairs of primers for qRT-PCR were designed using  
430 NCBI Primer-BLAST (<https://www.ncbi.nlm.nih.gov/tools/primer-blast>). In order to ensure  
431 the specificity of primers, UCSC InSilicon PCR (<http://genome.ucsc.edu/cgi-bin/hgPcr>) was  
432 used to compare the primer pairs with human genome (hg38). Some primer pairs were

433 specifically designed to span splicing sites (exon junctions). QRT-PCR were performed with  
434 SYBR Green master mix (Takara, Japan) on an ABI StepOnePlus (Applied Biosystems,  
435 USA). *GAPDH* as housekeeping gene was used as positive control. For each lncRNA, we  
436 selected one primer pair product of qRT-PCR for Sanger sequencing.

#### 437 **Quality control (QC) and normalization**

438 We calculated the read counts and transcripts per million (TPM) values using  
439 pseudoalignment of scRNA-seq reads to both protein-coding and lncRNA transcriptomes, as  
440 implemented in Kallisto (v0.46.0) [60] with default parameters, and summarized expression  
441 levels from the transcript level to the gene level.

442 Low-quality and doublet cells were removed if the number of expressed genes (counts of  
443 more than 1) was fewer than 2000 or higher than the medians of all cells plus  $3 \times$  the median  
444 absolute deviation, respectively. Moreover, the cells with the proportion of reads mapped to  
445 mitochondrial genes was larger than 10% were discarded. Genes with average counts of more  
446 than 1 and expressed in at least 1% of cells for each type of cancer were retained. The  
447 combined count tables from all T cells passing the above filtration were normalized using a  
448 pooling and deconvolution method implemented in the R package named  
449 *computeSumFactors* [61] with the sizes ranged from 80, 100, 120 to 140. According to the  
450 assumption that most genes were not differentially expressed, normalization was performed  
451 within each predefined cluster separately to compute cell size factors. The cell size factors  
452 were rescaled by normalization among clusters. Finally, the counts for each cell were  
453 normalized by dividing the cell counts by the cell size factor.

#### 454 **MetaCell modeling**

455 The MetaCell method [31], that partitioned the scRNA-seq dataset into disjointed and  
456 homogeneous cell groups (metacells) using the  $K$ -nn graph algorithm, was performed for  
457 both the CD8 and CD4 T cells independently. We first removed specific mitochondrial genes  
458 (annotated with the prefix “MT-”), that typically mark cells as being stressed or dying, rather  
459 than cellular identity. Based on the count matrices of both protein-coding and lncRNA genes,  
460 feature genes whose scaled variance (variance/mean on down-sampled matrices) exceeded  
461 0.08 were selected and used to compute cell-to-cell similarity using Pearson correlations.  
462 According to the cell-to-cell similarity matrices, two balanced  $K$ -nn similarity graphs for  
463 CD8 and CD4 T cells were constructed using the parameter  $K=100$  (the number of neighbors

464 for each cell was limited by  $K$ ). Next, we performed the resampling procedures (resampling  
465 75% of the cells in each iteration with 500 iterations) and co-clustering graph construction  
466 (the minimal cluster size was 50). Finally, the graphs of metacells (and the cells belonging to  
467 them) were projected into 2D spaces to explore the similarities between cells and metacells.

#### 468 **Annotation of metacells**

469 Annotation of metacells was performed based on the metacell confusion matrix and  
470 predefined cluster annotations (File S1) of T cells involved in the metacells. Briefly, we first  
471 created a hierarchical clustering of metacells according to the number of similarity  
472 relationships between their cells. Next, we generated clusters of metacells as confusion  
473 matrices based on the hierarchy results, then annotated these clusters according to the  
474 annotations of T cells.

#### 475 **Defining signature lncRNAs associated with T cell states**

476 To identify signature lncRNAs associated with effector and exhausted T cells as well as  
477 Tregs, as described in recent study [39], we adopted the anchor approach by identifying the  
478 lncRNAs that were significantly correlated to well-defined anchor genes, based on metacells'  
479 log enrichment scores (*lfp* values calculated by MetaCell method). The lncRNAs that  
480 significantly correlated with anchor genes (adjusted  $P$ -value  $< 0.01$  and ranked in the top 0.05  
481 percentile for each anchor gene) were regarded as signature lncRNAs. The anchor genes were  
482 defined as follows: the anchor genes of CD8 exhausted T cells included *HAVCR2*, *LAG3*,  
483 *PDCD1*, *TIGIT*, and *CTLA4*; the anchor genes of CD8 effector T cells included *CX3CRI*,  
484 *FGFBP2*, *GZMH* and *PRF1*; genes associated with Tregs included *FOXP3*; the anchor genes  
485 of CD4 exhausted T cells included *CXCL13*, *PDCD1*, *HAVCR2*, *TIGIT*, and *CTLA4*; genes  
486 associated with CD4 effector T cells included *GNLY*, *GZMB*, *GZMH*, *PRF1*, and *NKG7*.

#### 487 **Function prediction of signature lncRNAs based on co-expression network**

488 Based on *lfp* values of both lncRNA and protein-coding genes across all metacells, we used a  
489 custom pipeline for large-scale prediction of signature lncRNA functions by constructing the  
490 coding-lncRNA gene co-expression network [40, 41]. Briefly, genes with log enrichment  
491 scores ranked in the top 75% of each metacell were retained. Then,  $P$ -values of Pearson  
492 correlation coefficients for each gene pair were calculated based on the Fisher's asymptotic  
493 test using the *WGCNA* package of R.  $P$ -values were adjusted based on the Bonferroni  
494 multiple test correction using the *multtest* package of R. The gene pairs with an adjusted  $P$ -



495 value  $< 0.01$ , Pearson correlation coefficient  $> 0.7$ , and ranked in the top 5% for each gene  
496 were involved in co-expression network.

497 Based on the co-expression network, lncRNA functions were predicted using module-  
498 and hub-based methods. Specifically, the Markov cluster algorithm was adopted to identify  
499 co-expressed modules [40]. For each module, if the known genes were significantly enriched  
500 for at least one Gene Ontology (GO) term, the functions of the lncRNAs involved in the  
501 module were assigned as the same ones. For hub-based method, the functions of a hub  
502 lncRNA (node degree  $> 10$ ) were assigned, if its immediate neighboring genes were  
503 significantly enriched for at least one GO term.

#### 504 **Data availability**

505 All the novel lncRNA genes identified in current study and their expression files are available  
506 in the NONCODE database (<http://www.noncode.org/download.php>).

#### 507 **Authors' contributions**

508 HL, DB, JD, YZ and FL conceptualized and designed the study. HL and DB led the data  
509 analysis. HL performed the study and interpreted data. LJS performed experimental  
510 validation. YL collected the clinical samples and prepared the experimental materials. LS  
511 optimized the CNCI algorithm. WY, CW, XY and JW collected the data and performed T  
512 cell annotations. HL wrote the manuscript. HL, YZ and FL revised the manuscript. JD, YZ  
513 and FL supervised the project. All authors read and approved the final manuscript.

#### 514 **Competing interests**

515 The authors have declared no competing interests.

#### 516 **Acknowledgments**

517 This work was supported by the Science and Technology Project of Shenzhen (No.  
518 JHZ20170310090257380, JCYJ20170413092711058, JCYJ20170307095822325), China  
519 Postdoctoral Science Foundation (No. 2019M663369), Natural Science Foundation of  
520 Shenzhen (20190727160324164), and National Natural Science Foundation of China  
521 (31970636). We thank Dr. Lei Zhang and Yao He at Peking University for assistance with

522 raw data download. The data analysis was performed on AWS (China region) and we thank  
523 Mr Hansen Huang, Chao Wu and Fei Shi for providing the technical service and support.

#### 524 **Authors' ORCID IDs**

525 <sup>a</sup>0000-0003-3671-7786 (Luo, H)

526 <sup>b</sup>0000-0002-8833-5432 (Bu, D)

527 <sup>c</sup>0000-0001-7980-2003 (Shao, L)

528 <sup>d</sup>0000-0001-8199-8527 (Li, Y)

529 <sup>e</sup>0000-0002-5213-6941 (Sun, L)

530 <sup>f</sup>0000-0003-3048-7596 (Wang, C)

531 <sup>g</sup>0000-0001-7856-4533 (Wang, J)

532 <sup>h</sup>0000-0003-2138-5563 (Yang, W)

533 <sup>i</sup>0000-0001-6255-5181 (Yang, X)

534 <sup>j</sup>0000-0003-4064-3134 (Dong, J)

535 <sup>k</sup>0000-0001-6046-8420 (Zhao, Y)

536 <sup>l</sup>0000-0002-0606-8861 (Li, F)

#### 537 **References**

538

539 [1] Chen DS, Mellman I. Elements of cancer immunity and the cancer-immune set point.

540 Nature 2017;541:321-30.

541 [2] Bodor JN, Boumber Y, Borghaei H. Biomarkers for immune checkpoint inhibition in non-

542 small cell lung cancer (NSCLC). Cancer 2020;126:260-70.

543 [3] Houot R, Schultz LM, Marabelle A, Kohrt H. T-cell-based immunotherapy: adoptive cell

544 transfer and checkpoint inhibition. Cancer Immunol Res 2015;3:1115-22.

- 545 [4] Restifo NP, Dudley ME, Rosenberg SA. Adoptive immunotherapy for cancer: harnessing  
546 the T cell response. *Nat Rev Immunol* 2012;12:269-81.
- 547 [5] Hilmi M, Vienot A, Rousseau B, Neuzillet C. Immune therapy for liver cancers. *Cancers*  
548 (*Basel*) 2019;12.
- 549 [6] Carter JA, Gilbo P, Atwal GS. IMPRES does not reproducibly predict response to  
550 immune checkpoint blockade therapy in metastatic melanoma. *Nat Med* 2019;25:1833-5.
- 551 [7] Hugo W, Zaretsky JM, Sun L, Song C, Moreno BH, Hu-Lieskovan S, et al. Genomic and  
552 transcriptomic features of response to anti-*PD-1* therapy in metastatic melanoma. *Cell*  
553 2016;165:35-44.
- 554 [8] Auslander N, Zhang G, Lee JS, Frederick DT, Miao B, Moll T, et al. Robust prediction of  
555 response to immune checkpoint blockade therapy in metastatic melanoma. *Nat Med*  
556 2018;24:1545-9.
- 557 [9] Rinn JL, Chang HY. Genome regulation by long noncoding RNAs. *Annu Rev Biochem*  
558 2012;81:145-66.
- 559 [10] Fang S, Zhang L, Guo J, Niu Y, Wu Y, Li H, et al. NONCODEV5: a comprehensive  
560 annotation database for long non-coding RNAs. *Nucleic Acids Res* 2018;46:D308-D14.
- 561 [11] Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al.  
562 GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res*  
563 2012;22:1760-74.
- 564 [12] Iyer MK, Niknafs YS, Malik R, Singhal U, Sahu A, Hosono Y, et al. The landscape of  
565 long noncoding RNAs in the human transcriptome. *Nat Genet* 2015;47:199-208.
- 566 [13] Yu WD, Wang H, He QF, Xu Y, Wang XC. Long noncoding RNAs in cancer-immunity  
567 cycle. *J Cell Physiol* 2018;233:6518-23.
- 568 [14] Wang L, Felts SJ, Van Keulen VP, Scheid AD, Block MS, Markovic SN, et al.  
569 Integrative genome-wide analysis of long noncoding RNAs in diverse immune cell types of  
570 melanoma patients. *Cancer Res* 2018;78:4411-23.
- 571 [15] Xu J, Cao X. Long noncoding RNAs in the metabolic control of inflammation and  
572 immune disorders. *Cell Mol Immunol* 2019;16:1-5.
- 573 [16] Spurlock CF, Crooke PS, Aune TM. Biogenesis and transcriptional regulation of long  
574 noncoding RNAs in the human immune system. *J Immunol* 2016;197:4509-17.
- 575 [17] Fanucchi S, Fok ET, Dalla E, Shibayama Y, Borner K, Chang EY, et al. Immune genes  
576 are primed for robust transcription by proximal long noncoding RNAs located in nuclear  
577 compartments. *Nat Genet* 2019;51:138-50.

- 578 [18] Huang D, Chen J, Yang L, Ouyang Q, Li J, Lao L, et al. NKILA lncRNA promotes  
579 tumor immune evasion by sensitizing T cells to activation-induced cell death. *Nat Immunol*  
580 2018;19:1112-25.
- 581 [19] Gomez JA, Wapinski OL, Yang YW, Bureau JF, Gopinath S, Monack DM, et al. The  
582 NeST long ncRNA controls microbial susceptibility and epigenetic activation of the  
583 interferon-gamma locus. *Cell* 2013;152:743-54.
- 584 [20] Ren X, Kang B, Zhang Z. Understanding tumor ecosystems by single-cell sequencing:  
585 promises and limitations. *Genome Biol* 2018;19:211.
- 586 [21] Chung W, Eum HH, Lee HO, Lee KM, Lee HB, Kim KT, et al. Single-cell RNA-seq  
587 enables comprehensive tumour and immune cell profiling in primary breast cancer. *Nat*  
588 *Commun* 2017;8:15081.
- 589 [22] Lavin Y, Kobayashi S, Leader A, Amir ED, Elefant N, Bigenwald C, et al. Innate  
590 immune landscape in early lung adenocarcinoma by paired single-cell analyses. *Cell*  
591 2017;169:750-65 e17.
- 592 [23] Zhang Q, He Y, Luo N, Patel SJ, Han Y, Gao R, et al. Landscape and dynamics of single  
593 immune cells in hepatocellular carcinoma. *Cell* 2019;179:829-45 e20.
- 594 [24] Savas P, Virassamy B, Ye C, Salim A, Mintoff CP, Caramia F, et al. Single-cell  
595 profiling of breast cancer T cells reveals a tissue-resident memory subset associated with  
596 improved prognosis. *Nat Med* 2018;24:986-93.
- 597 [25] Azizi E, Carr AJ, Plitas G, Cornish AE, Konopacki C, Prabhakaran S, et al. Single-cell  
598 map of diverse immune phenotypes in the breast tumor microenvironment. *Cell*  
599 2018;174:1293-308 e36.
- 600 [26] Uszczyńska-Ratajczak B, Lagarde J, Frankish A, Guigo R, Johnson R. Towards a  
601 complete map of the human long non-coding RNA transcriptome. *Nat Rev Genet*  
602 2018;19:535-48.
- 603 [27] Stegle O, Teichmann SA, Marioni JC. Computational and analytical challenges in  
604 single-cell transcriptomics. *Nat Rev Genet* 2015;16:133-45.
- 605 [28] Vallejos CA, Risso D, Scialdone A, Dudoit S, Marioni JC. Normalizing single-cell RNA  
606 sequencing data: challenges and opportunities. *Nat Methods* 2017;14:565-71.
- 607 [29] Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively  
608 parallel digital transcriptional profiling of single cells. *Nat Commun* 2017;8:14049.
- 609 [30] Picelli S, Faridani OR, Bjorklund AK, Winberg G, Sagasser S, Sandberg R. Full-length  
610 RNA-seq from single cells using Smart-seq2. *Nat Protoc* 2014;9:171-81.

- 611 [31] Baran Y, Bercovich A, Sebe-Pedros A, Lubling Y, Giladi A, Chomsky E, et al. MetaCell:  
612 analysis of single-cell RNA-seq data using K-nn graph partitions. *Genome Biol* 2019;20:206.
- 613 [32] Guo X, Zhang Y, Zheng L, Zheng C, Song J, Zhang Q, et al. Global characterization of  
614 T cells in non-small-cell lung cancer by single-cell sequencing. *Nat Med* 2018;24:978-85.
- 615 [33] Zhang L, Yu X, Zheng L, Zhang Y, Li Y, Fang Q, et al. Lineage tracking reveals  
616 dynamic relationships of T cells in colorectal cancer. *Nature* 2018;564:268-72.
- 617 [34] Zheng C, Zheng L, Yoo JK, Guo H, Zhang Y, Guo X, et al. Landscape of infiltrating T  
618 cells in liver cancer revealed by single-cell sequencing. *Cell* 2017;169:1342-56 e16.
- 619 [35] Perteau M, Perteau GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. StringTie  
620 enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*  
621 2015;33:290-5.
- 622 [36] Kong L, Zhang Y, Ye ZQ, Liu XQ, Zhao SQ, Wei L, et al. CPC: assess the protein-  
623 coding potential of transcripts using sequence features and support vector machine. *Nucleic*  
624 *Acids Res* 2007;35:W345-9.
- 625 [37] Sun L, Luo H, Bu D, Zhao G, Yu K, Zhang C, et al. Utilizing sequence intrinsic  
626 composition to classify protein-coding and long non-coding transcripts. *Nucleic Acids Res*  
627 2013;41:e166.
- 628 [38] Jiang S, Cheng SJ, Ren LC, Wang Q, Kang YJ, Ding Y, et al. An expanded landscape of  
629 human long noncoding RNA. *Nucleic Acids Res* 2019;47:7842-56.
- 630 [39] Li H, van der Leun AM, Yofe I, Lubling Y, Gelbard-Solodkin D, van Akkooi ACJ, et al.  
631 Dysfunctional CD8 T cells form a proliferative, dynamically regulated compartment within  
632 human melanoma. *Cell* 2019;176:775-89 e18.
- 633 [40] Liao Q, Liu C, Yuan X, Kang S, Miao R, Xiao H, et al. Large-scale prediction of long  
634 non-coding RNA functions in a coding-non-coding gene co-expression network. *Nucleic*  
635 *Acids Res* 2011;39:3864-78.
- 636 [41] Luo H, Bu D, Sun L, Fang S, Liu Z, Zhao Y. Identification and function annotation of  
637 long intervening noncoding RNAs. *Brief Bioinform* 2017;18:789-97.
- 638 [42] Wherry EJ, Kurachi M. Molecular and cellular insights into T cell exhaustion. *Nat Rev*  
639 *Immunol* 2015;15:486-99.
- 640 [43] Wright MD, Ni J, Rudy GB. The *L6* membrane proteins--a new four-transmembrane  
641 superfamily. *Protein Sci* 2000;9:1594-600.
- 642 [44] Chang YW, Chen SC, Cheng EC, Ko YP, Lin YC, Kao YR, et al. *CD13*  
643 (aminopeptidase N) can associate with tumor-associated antigen *L6* and enhance the motility  
644 of human lung cancer cells. *Int J Cancer* 2005;116:243-52.

- 645 [45] Lekishvili T, Fromm E, Mujoomdar M, Berditchevski F. The tumour-associated antigen  
646 *L6 (L6-Ag)* is recruited to the tetraspanin-enriched microdomains: implication for tumour cell  
647 motility. *J Cell Sci* 2008;121:685-94.
- 648 [46] Allioli N, Vincent S, Vlaeminck-Guillem V, Decaussin-Petrucci M, Ragage F, Ruffion  
649 A, et al. *TM4SF1*, a novel primary androgen receptor target gene over-expressed in human  
650 prostate cancer and involved in cell migration. *Prostate* 2011;71:1239-50.
- 651 [47] Yang W, Bai Y, Xiong Y, Zhang J, Chen S, Zheng X, et al. Potentiating the antitumour  
652 response of CD8(+) T cells by modulating cholesterol metabolism. *Nature* 2016;531:651-5.
- 653 [48] Marton N, Baricza E, Ersek B, Buzas EI, Nagy G. The Emerging and diverse roles of  
654 *Src*-like adaptor proteins in health and disease. *Mediators Inflamm* 2015;2015:952536.
- 655 [49] Sosinowski T, Pandey A, Dixit VM, Weiss A. *Src*-like adaptor protein (*SLAP*) is a  
656 negative regulator of T cell receptor signaling. *J Exp Med* 2000;191:463-74.
- 657 [50] Pandey A, Ibarrola N, Kratchmarova I, Fernandez MM, Constantinescu SN, Ohara O, et  
658 al. A novel *Src* homology 2 domain-containing molecule, *Src*-like adapter protein-2 (*SLAP-2*),  
659 which negatively regulates T cell receptor signaling. *J Biol Chem* 2002;277:19131-8.
- 660 [51] Park SK, Beaven MA. Mechanism of upregulation of the inhibitory regulator, *src*-like  
661 adaptor protein (*SLAP*), by glucocorticoids in mast cells. *Mol Immunol* 2009;46:492-7.
- 662 [52] Yao RW, Wang Y, Chen LL. Cellular functions of long noncoding RNAs. *Nat Cell Biol*  
663 2019;21:542-51.
- 664 [53] Agirre X, Meydan C, Jiang Y, Garate L, Doane AS, Li Z, et al. Long non-coding RNAs  
665 discriminate the stages and gene regulatory states of human humoral immune response. *Nat*  
666 *Commun* 2019;10:821.
- 667 [54] Liu F, Zhang Y, Zhang L, Li Z, Fang Q, Gao R, et al. Systematic comparative analysis  
668 of single-nucleotide variant detection methods from single-cell RNA sequencing data.  
669 *Genome Biol* 2019;20:242.
- 670 [55] Quek XC, Thomson DW, Maag JL, Bartonicek N, Signal B, Clark MB, et al. IncRNADB  
671 v2.0: expanding the reference database for functional long noncoding RNAs. *Nucleic Acids*  
672 *Res* 2015;43:D168-73.
- 673 [56] Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast  
674 universal RNA-seq aligner. *Bioinformatics* 2013;29:15-21.
- 675 [57] Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence  
676 alignment/map format and SAMtools. *Bioinformatics* 2009;25:2078-9.

677 [58] Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al.  
678 Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and  
679 isoform switching during cell differentiation. *Nat Biotechnol* 2010;28:511-5.

680 [59] Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, et al.  
681 GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res*  
682 2019;47:D766-D73.

683 [60] Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq  
684 quantification. *Nat Biotechnol* 2016;34:525-7.

685 [61] Lun AT, Bach K, Marioni JC. Pooling across cells to normalize single-cell RNA  
686 sequencing data with many zero counts. *Genome Biol* 2016;17:75.

687

## 688 **Figure Legends**

689 **Figure 1 The statistics of assembled transcripts and workflow for novel lncRNA**  
690 **identification process in T cells during cancer immunity**

691 **A.** Violin plots showing the number of assembled transcripts and the number of those  
692 matched to the reference at single cell level across five HCC patients. **B.** Number of  
693 assembled transcripts that matched to reference across five HCC patients based on four  
694 different strategies. \*\*\* indicates  $P$ -value  $< 0.001$  (Wilcoxon rank sum test). **C.** Correlation  
695 of the number of cells and the number of assembled transcripts across different subsets for  
696 CRC, HCC and NSCLC. A 95% confidence interval was added and shown as coloured  
697 regions. **D.** Scheme of pipeline used to identify the novel lncRNAs expressed in T cells  
698 during cancer immunity using three full-length scRNA-seq datasets. **E.** The statistics of  
699 assembled transcripts that matched to reference protein-coding and reference lncRNA genes.  
700 CRC, colorectal cancer; HCC, hepatocellular carcinoma; NSCLC, non-small-cell lung cancer;  
701 P, peripheral blood; N, adjacent normal tissue; T, tumor tissue.

702 **Figure 2 Single T cell sorting and quality evaluation of an example novel lncRNA**

703 **A.** The results of flow cytometric analysis. CD8 and CD4 T cells from three patients were  
704 separated by magnetic beads and stained with flow cytometry antibody CD8-APC and CD4-  
705 APC respectively (Isotype was used as negative control). **B.** An example of novel intergenic  
706 lncRNA that was validated by Sanger sequencing. The genomic views are generated from  
707 UCSC genome browser. The spliced sequence outputted by Sanger sequencing is shown.

708 **Figure 3 Characterization of lncRNA expression patterns at single-cell level**

709 **A.** The number of protein-coding, reference lncRNA, and novel lncRNA genes expressed in  
710 T cells across three cancer types. \*\*\* indicates  $P$ -value  $< 0.001$  (Wilcoxon rank sum test). **B.**  
711 The plots show the percentage of expressing cells against the mean expression level  
712 (logCounts) for protein-coding, reference lncRNA, and novel lncRNA genes across three  
713 cancer types. The numbers of genes that are expressed in at least 25% of cells are labelled.

714 **Figure 4 Characterization of T cell states based on 2D projection of T cells and the**  
715 **annotation of metacell maps**

716 **A.** 2D projection of CD8 T cells from three cancer types into 43 metacells. **B.** 2D projection  
717 of CD4 T cells from three cancer types into 65 metacells. **C, D.** CD8 (**C**) and CD4 (**D**)  
718 metacells (rows) are ordered by groups and organized within each group. The first panel of  
719 the bar plot shows the number of cells of different clusters in each metacell. The second and  
720 third panel of the bar plots show the percentage of cells from different cancer types and  
721 tissues in each metacell respectively. Heatmaps show the confusion matrix (the pairwise  
722 similarities between metacells) for CD8 (**C**) and CD4 (**D**) metacells. The annotations of  
723 different metacell groups are shown on the right. **E, F.** 2D projections of the composition of  
724 CD8 (**E**) and CD4 (**F**) T cells from different clusters. P, peripheral blood; N, adjacent normal  
725 tissue; T, tumor tissue.

726 **Figure 5 The correlation and expression analyses of signature lncRNAs associated with**  
727 **different T cell states**

728 **A, B.** Gene-gene correlation heatmap for signature lncRNA and anchor genes in CD8 (**A**) and  
729 CD4 (**B**) T cells. The signature gene modules and two anchor genes (*CTLA4* and *FOXP3*) are  
730 labelled on the right. **C.** Expression of signature lncRNA and anchor genes across CD8  
731 metacells. Metacells and metacell groups associated with effector and exhausted functions are  
732 shown on the bottom. The anchor genes are marked with red color on the right.

733 **Figure 6 Functional annotation analyses of signature lncRNAs**

734 **A-C.** Functional enrichment maps of CD8 effector/exhausted (**A**), CD4 effector/exhausted (**B**)  
735 and CD4 Treg (**C**) signature lncRNAs. The enriched gene sets from Gene Ontology based on  
736 the predicted functions of signature lncRNA genes are visualized by Cytoscape plugin  
737 Enrichment Map. Each node represents a gene set; size of the node is indicative of the



738 number of genes and the color intensity reflects the level of significance. Effector signature  
739 gene sets are shown in red circles, exhausted or Treg ones in green and the common gene sets  
740 in orange. Maps are differently magnified for easier visualization. **D-F**. The genomic view  
741 (**D**), co-expressed genes (**E**) and functional annotations (**F**) of effector signature lncRNA  
742 *TM4SF19-AS1*. **G-I**. The genomic view (**G**), co-expressed genes (**H**) and functional  
743 annotations (**I**) of exhausted signature lncRNA *XLOC-633950* (novel). The genomic views  
744 are generated from UCSC genome browser. In (**E**) and (**H**), co-expressed protein-coding,  
745 reference lncRNA and novel lncRNA genes are colored by pink, light green and light yellow  
746 respectively.

## 747 **Supplementary material**

### 748 **Figure S1 The statistics of T cell data analysis**

749 **A**. The number of cells in different subsets across all patients from three cancer types. **B, C**.  
750 The number (**B**) and the ratio (**C**) of uniquely mapped read pairs of T cell sequencing data. **D**.  
751 The number of splices of mapping results. **E**. The different strategies used to explore the best  
752 way to obtain novel transcripts. **F**. The number of assembled transcripts in each subset. PTC,  
753 CD8<sup>+</sup> cytotoxic T cells from peripheral blood; TTC, CD8<sup>+</sup> cytotoxic T cells from tumor  
754 tissue; NTC, CD8<sup>+</sup> cytotoxic T cells from adjacent normal tissue; PTH, CD4<sup>+</sup>CD25<sup>-</sup> cells  
755 from peripheral blood; TTH, CD4<sup>+</sup>CD25<sup>-</sup> cells from tumor tissue; NTH, CD4<sup>+</sup>CD25<sup>-</sup> cells  
756 from adjacent normal tissue; PTR, CD4<sup>+</sup>CD25<sup>hi</sup> cells from peripheral blood; TTR,  
757 CD4<sup>+</sup>CD25<sup>hi</sup> cells from tumor tissue; NTR, CD4<sup>+</sup>CD25<sup>hi</sup> cells from adjacent normal tissue;  
758 PTY, CD4<sup>+</sup>CD25<sup>int</sup> cells from peripheral blood; TTY, CD4<sup>+</sup>CD25<sup>int</sup> cells from tumor tissue;  
759 NTY, CD4<sup>+</sup>CD25<sup>int</sup> cells from adjacent normal tissue; PPQ, CD4<sup>+</sup> T cells from peripheral  
760 blood; TPQ, CD4<sup>+</sup> T cells from tumor tissue; NPQ, CD4<sup>+</sup> T cells from adjacent normal tissue;  
761 CRC, colorectal cancer; HCC, hepatocellular carcinoma; NSCLC, non-small-cell lung cancer.

### 762 **Figure S2 The cluster hierarchy of metacells**

763 **A, B**. The cluster hierarchy of CD8 (**A**) and CD4 (**B**) metacells. Subtrees in blue, sibling  
764 subtrees in gray. The metacells are colored and labelled on bottom.

### 765 **Figure S3 2D projections of CD8 T cells**

766 **A, B.** The composition of CD8 T cells from different clusters (**A**) and cancer types (**B**).  
767 Metacells and the cells involved in them are marked by different colors. The number of cells  
768 within each cluster is shown in brackets.

769 **Figure S4 2D projections of CD4 T cells**

770 **A, B.** The composition of CD4 T cells from different clusters (**A**) and cancer types (**B**).  
771 Metacells and the cells involved in them are marked by different colors. The number of cells  
772 within each cluster is shown in brackets.

773 **Figure S5 Expression of signature lncRNA and anchor genes across CD4 metacells**

774 Metacells and metacell groups associated with effector, exhausted and Treg functions are  
775 shown on the bottom. The anchor genes are marked with red color on the right.

776 **Figure S6 Functional enrichment maps of shared signature lncRNAs**

777 **A-C.** Functional enrichment maps of shared signature lncRNAs between CD8 effector and  
778 CD4 effector function (**A**), between CD8 exhausted and CD4 exhausted function (**B**) and  
779 between CD4 exhausted and CD4 Treg function (**C**). Each node represents a gene set; size of  
780 the node is indicative of the number of genes and the color intensity reflects the level of  
781 significance. Maps are differently magnified for easier visualization.

782 **Table S1 The basic information of single T cell data**

783 **Table S2 Clinical characteristics of three cancer patients**

784 **Table S3 The list of novel lncRNAs successfully validated by Sanger sequencing**

785 **Table S4 The list of specific-expressed lncRNAs**

786 **Table S5 The composition of CD8 metacells**

787 **Table S6 The composition of CD4 metacells**

788 **Table S7 The list of signature lncRNAs**

789 **Table S8 Functional annotations of 84 signature lncRNAs**

790 **Table S9 Functional enrichment results of CD8 effector/exhausted signature lncRNAs**

791 **Table S10 Functional enrichment results of CD4 effector/exhausted signature lncRNAs**

792 **Table S11 Functional enrichment results of CD4 Treg signature lncRNAs**

793 Supplementary Table1-8 are Excel format, and Supplementary Table9-11 are Word format.













